

1

2 Review and further developments in statistical corrections for

3 Winner's Curse in genetic association studies

4

5

6 Amanda Forde<sup>1\*</sup>, Gibran Hemani<sup>2,3</sup>, John Ferguson<sup>4</sup>

7

8 <sup>1</sup>School of Mathematical and Statistical Sciences, University of Galway, Galway, Ireland

9 <sup>2</sup> MRC Integrative Epidemiology Unit, University of Bristol, Oakfield House, Oakfield

10 Grove, Bristol, BS8 2BN, UK

11 <sup>3</sup> Population Health Sciences, Bristol Medical School, University of Bristol, Bristol, UK

12 <sup>4</sup>HRB Clinical Research Facility, University of Galway, Galway, Ireland

13

14

15 \* Corresponding author

16 E-mail: [a.forde21@nuigalway.ie](mailto:a.forde21@nuigalway.ie)

## 1 **Abstract**

2           Genome-wide association studies (GWAS) are commonly used to identify genomic  
3 variants that are associated with complex traits, and estimate the magnitude of this  
4 association for each variant. However, it has been widely observed that the association  
5 estimates of variants tend to be lower in a replication study than in the study that discovered  
6 those associations. A phenomenon known as *Winner's Curse* is responsible for this upward  
7 bias present in association estimates of significant variants in the discovery study. We review  
8 existing *Winner's Curse* correction methods which require only GWAS summary statistics in  
9 order to make adjustments. In addition, we propose modifications to improve existing  
10 methods and propose a novel approach which uses the parametric bootstrap. We evaluate and  
11 compare methods, first using a wide variety of simulated data sets and then, using real data  
12 sets for three different traits. The metric, estimated mean squared error (MSE) over  
13 significant SNPs, was primarily used for method assessment. Our results indicate that widely  
14 used conditional likelihood based methods tend to perform poorly. The other considered  
15 methods behave much more similarly, with our proposed bootstrap method demonstrating  
16 very competitive performance. To complement this review, we have developed an R package,  
17 'winnerscurse' which can be used to implement these various *Winner's Curse* adjustment  
18 methods to GWAS summary statistics.

## 19 **Author Summary**

20           A genome-wide association study is designed to analyse many common genetic  
21 variants in thousands of samples and identify which variants are associated with a trait of  
22 interest. It provides estimates of association strength for each variant and variants are  
23 classified as associated if their test statistics obtained in the study pass a chosen significance

24 threshold. However, due to a phenomenon known as *Winner's Curse*, the association  
25 estimates of these significant variants tend to be upward biased and greater in magnitude than  
26 their true values. Naturally, this bias has adverse consequences for downstream statistical  
27 techniques which use these estimates. In this paper, we look at current methods which have  
28 been designed to combat *Winner's Curse* and propose modifications to these methods in  
29 order to improve performance. Using a wide variety of simulated data sets as well as real  
30 data, we perform a thorough evaluation of these methods. We use a metric which allows us to  
31 identify which methods, on average, produce adjusted estimates for significant variants that  
32 are closest to the true values. To accompany our work, we have created an R package,  
33 'winnerscurse', which allows users to easily apply *Winner's Curse* correction methods to  
34 their data sets.

## 35 **Introduction**

36 It has been observed that in general, the effect size of a variant or single nucleotide  
37 polymorphism (SNP) tends to be lower in a replication study than in the genome-wide  
38 association study (GWAS) that discovered the SNP-trait association. This observation is due  
39 to the phenomenon known as *Winner's Curse*. In the context of a single discovery GWAS,  
40 the term *Winner's Curse* describes how the estimates of association strength for SNPs that  
41 have been deemed most significant are very likely to be exaggerated compared with their true  
42 underlying values. These estimated effect sizes can take the form of log odds ratios (log-OR)  
43 resulting from a logistic regression for a binary outcome, e.g. disease status, or regression  
44 coefficients (beta) derived from a linear regression for a quantitative trait.

45 Dudbridge & Newcombe (1) detail two sources of *Winner's Curse* in GWASs,  
46 namely ranking bias and selection bias. Ranking bias stems from ranking many SNPs, often  
47 close to a million or more, by some measure of effect size or statistical significance. In

48 practice,  $p$ -values are generally used. It is then expected that the bias will be greatest for those  
49 variants which have been ranked highly. Selection bias describes how the use of a stringent  
50 threshold, such as  $5 \times 10^{-8}$ , can result in overestimated effect sizes for SNPs that exceed this  
51 threshold.

52 *Winner's Curse* bias can have many practical consequences, especially with respect to  
53 techniques which are reliant on SNP-trait association estimates obtained from GWASs. One  
54 such example is Mendelian randomization (MR), a statistical framework which uses genetic  
55 variants as instrumental variables to estimate the magnitude of the casual effect of an  
56 exposure on an outcome. In the case of two-sample MR, if the same GWAS is used to  
57 identify instrument SNPs and estimate their effects relative to the exposure, *Winner's Curse*  
58 will result in the over-estimation of these SNP-exposure associations. This bias will then  
59 propagate into the causal estimate, resulting in a deflation of this estimate. On the other hand,  
60 if instrument SNPs are discovered in the same GWAS as that used to estimate the SNP-  
61 outcome associations, the causal estimate will be inflated (2). In addition, *Winner's Curse* has  
62 been shown to greatly increase the magnitude of weak instrument bias in these MR analyses  
63 (3). Another implication of *Winner's Curse* bias is in the use of polygenic risk scores which  
64 employs GWAS results for prediction purposes. Enlarged association estimates of significant  
65 variants used in creating the polygenic score can lead to reduced accuracy in out-of-sample  
66 prediction (4).

67 In this paper, we review existing *Winner's Curse* correction methods and explore  
68 possible modifications that could be made in order to improve these methods. However,  
69 eliminating this bias induced by *Winner's Curse* is known to be a difficult task. Several bias  
70 reduction approaches have been proposed in recent years, with one of the earliest being the  
71 Conditional Likelihood method suggested by Ghosh et al. (5). This method makes an  
72 adjustment to the association estimate of each SNP which has been deemed significant, i.e.

73 those with  $p$ -values less than the specified genome-wide significance threshold. In contrast to  
74 this approach in which the correction is performed to each SNP separately, independently of  
75 estimated associations of other SNPs, alternative methods have been suggested which involve  
76 the use of all SNPs, including those which do not pass the threshold, in order to produce bias-  
77 reduced estimated effect sizes. The empirical Bayes method described by Ferguson et al. (6)  
78 determines a suitable correction for each SNP by using the collective distribution of all effect  
79 sizes. Bigdeli et al. (7) suggested the use of FDR Inverse Quantile Transformation (FIQT),  
80 while Faye et al. (8) proposed a bootstrap shrinkage estimator with application to the GWAS  
81 setting. As this bootstrap approach requires individual-level data, we propose an alternative  
82 form of this method which uses bootstrapping with summary statistics to make corrections.

83         The focus in this paper is on methods which attempt to reduce the effect of the bias  
84 induced by *Winner's Curse* using only GWAS summary statistics, not individual-level data,  
85 the reason being that approaches based on summary data tend to be more computationally  
86 efficient, in terms of run time and memory efficiency. Furthermore, GWAS summary  
87 statistics are much more accessible and are more widely used in epidemiological techniques  
88 such as MR. In addition, there exist methods which use both a discovery and a replication  
89 GWAS in order to make suitable corrections to estimated effect sizes of significant SNPs.  
90 Examples include the UMVCUE of Bowden and Dudbridge (9) and an additional conditional  
91 likelihood method, Zhong and Prentice (10). That said, the concentration of our work detailed  
92 here is on techniques which have been designed for use when a replication sample is  
93 unavailable.

94         As mentioned above, we have made amendments to existing *Winner's Curse*  
95 correction methods to address certain weaknesses. In particular, we investigated  
96 modifications that could be made to the empirical Bayes method in order to ensure that it  
97 makes better adjustments to association estimates. Following this review of correction

98 methods, a rigorous evaluation and comparison of these methods was performed. This  
99 assessment took place by means of a simulation study as well as engagement with three real  
100 data sets. Simulations allowed us to compare methods easily over a wide range of different  
101 possible genetic architectures. We then used UK Biobank (UKBB) body mass index (BMI),  
102 type 2 diabetes (T2D) and height data sets to see how these techniques would perform in  
103 more realistic settings in which a large degree of linkage disequilibrium (LD) exists. In both  
104 instances, assessment of methods was predominantly based on the computation of estimated  
105 mean squared error (MSE) over significant SNPs. A notable challenge that was encountered  
106 at the start of the work discussed in this paper was the lack of available software to  
107 implement these various correction methods. Therefore, to complement this review, we have  
108 developed an R package, namely ‘winnerscurse’  
109 (<https://github.com/amandaforde/winnerscurse>), which can be used to apply a number of  
110 *Winner’s Curse* adjustment methods to GWAS summary statistics. Techniques which require  
111 a replication GWAS are also included in this package.

## 112 **Materials and methods**

113 Throughout this paper, we let  $Z_i = \frac{\hat{\beta}_i}{\widehat{se}(\hat{\beta}_i)}$  and  $\mu_i = \frac{\beta_i}{\widehat{se}(\hat{\beta}_i)}$  with the assumption that

$$Z_i \sim N(\mu_i, 1) \quad (1)$$

114 asymptotically, in which  $\beta_i$  denotes the true effect size of SNP  $i$  for  $i = 1, \dots, N$ ,  $\hat{\beta}_i$  its  
115 estimated effect size, with respect to a trait of interest, and  $\widehat{se}(\hat{\beta}_i)$  its estimated standard error.  
116  $N$  represents the total number of SNPs in the discovery GWAS. Depending on the type of  
117 phenotype, be it a disease or a quantitative trait, this estimated effect size can represent a log  
118 odds ratio or a regression coefficient attained from a linear regression, respectively.

## 119 **Conditional Likelihood**

120 As mentioned, the conditional likelihood method of Ghosh et al. (5) notably differs in  
 121 its approach at making *Winner's Curse* corrections from the other methods evaluated in this  
 122 paper. Adjustments are made to the estimated effect sizes of only those SNPs which satisfy  $|z|$   
 123  $> c$ , where  $c$  is the value corresponding to the pre-specified significance threshold. The  
 124 reduction in estimated effect size for each significant SNP is imposed independently of other  
 125 SNPs and is directly determined by the value of  $c$ . Recognizing that a SNP has been deemed  
 126 significant, the corresponding conditional likelihood is given by:

$$L_c(\mu_i) = p_{\mu_i}(z_i | |Z_i| > c) = \frac{p_{\mu_i}(z_i)}{P_{\mu_i}(|Z_i| > c)} = \frac{\phi(z_i - \mu_i)}{\Phi(-c + \mu_i) + \Phi(-c - \mu_i)} \quad (2)$$

127 in which  $\phi(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$ , the probability density function of the standard Gaussian  
 128 distribution and  $\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-t^2/2} dt$ , the corresponding cumulative distribution  
 129 function (cdf). In general,  $c$  takes the form of  $c = \Phi^{-1}\left(1 - \frac{\alpha}{2}\right)$ , with  $\alpha$  being a threshold to  
 130 which a Bonferroni correction has been applied in order to control for family-wise error rate,  
 131 e.g.  $5 \times 10^{-8}$ .

132 Using this conditional likelihood, three estimators of  $\mu_i$ , or equivalently of  $\beta_i$  as any  
 133 estimator for  $\mu_i$  can be used to produce an estimator for  $\beta_i$  by simply multiplying it by  $\widehat{se(\beta_i)}$ ,  
 134 were proposed. The first,  $\tilde{\mu}_{i1}$  is the obvious conditional maximum likelihood estimator:

$$\tilde{\mu}_{i1} = \arg \max_{\mu_i} L_c(\mu_i) \quad (3)$$

135 while the second,  $\tilde{\mu}_{i2}$  is defined as

$$\tilde{\mu}_{i2} = \frac{\int_{-\infty}^{\infty} \mu_i L_c(\mu_i) d\mu_i}{\int_{-\infty}^{\infty} L_c(\mu_i) d\mu_i}. \quad (4)$$

136 This is the mean of the random variable that follows the distribution  $L_c(\mu_i)$ , normalized to  
 137 ensure a proper density. However, it was observed that for instances in which the true effect  
 138 size,  $\beta_i$  is close to that of a null effect, the estimator  $\tilde{\mu}_{i2}$  has greater mean squared error than

139  $\tilde{\mu}_{i1}$  but for true effect sizes further from zero,  $\tilde{\mu}_{i2}$  performs better. Therefore, the use of  $\tilde{\mu}_{i3} =$   
140  $\frac{\tilde{\mu}_{i1} + \tilde{\mu}_{i2}}{2}$ , which can combine the strengths of these two estimators in order to curtail *Winner's*  
141 *Curse* bias for significant SNPs more accurately, was suggested.

## 142 **Empirical Bayes**

143 Motivated by Efron's empirical Bayes implementation of Tweedie's formula to  
144 correct for selection bias (11), the empirical Bayes method detailed by Ferguson et al. (6)  
145 focuses on the importance of sharing information between SNPs in order to make  
146 adjustments, through the exploitation of the empirical distribution of all effect sizes. This is a  
147 notably different approach to that of the previously discussed conditional likelihood method,  
148 which when making a correction to the estimated effect size of a particular SNP essentially  
149 fails to acknowledge the existence of any other SNPs.

150 Under the normal sampling assumption described by Eq (16), Tweedie's formula  
151 describes the relationship between the posterior mean,  $E(\mu|z)$ , and the marginal density  
152 function,  $p(z)$ , as

$$E(\mu|z) = z + \frac{d}{dz} \log p(z) \quad (5)$$

153 Amazingly, provided one can estimate  $p(z)$ , Tweedie's formula facilitates estimation of the  
154 posterior mean in complete absence of knowledge of the prior distribution,  $p(\mu)$ , which in this  
155 instance is the true distribution of standardized effect sizes across the genome. Thus, the  
156 estimator of  $\mu_i$  proposed by this method takes the form of

$$\tilde{\mu}_i = E(\widehat{\mu}_i|z_i) = z_i + \frac{d}{dz_i} \log \widehat{p}(z_i). \quad (6)$$

157 Estimation of  $\log p(z)$  occurs upon application of the following steps. First, partitions  
158 of the interval  $[z_1, z_N]$  of identical width are formed, in which the  $z$ -statistics have been



159 arranged in ascending order. The number of  $z$ -statistics which fall inside each partition are  
160 noted and regressed against a set of natural cubic spline basis functions with knots located at  
161 the midpoint of each partition, using a Poisson generalized linear model. Ferguson et al. (6)  
162 suggest choosing the number of basis functions so that the Bayesian Information Criterion  
163 (BIC) is minimized for the model. The fitted regression function at  $z$  is then used to obtain the  
164 estimate for  $\log p(z)$ , and subsequently,  $\tilde{\mu}_i$  for  $i = 1, \dots, N$ , by means of numerical  
165 differentiation.

166 Ferguson et al. (6) show that if the true marginal density,  $p(z)$ , could be used here,  
167 then the empirical Bayes estimator would perform optimally at minimizing the mean squared  
168 error (MSE) over all SNPs. However, since it is only an estimate of  $p(z)$  that can be obtained,  
169 this optimal behaviour is not guaranteed. This is especially a concern in the extreme tails of  
170 the distribution where the  $z$ -statistics of the most significant SNPs lie as it is more difficult to  
171 accurately estimate  $p(z)$  in these regions. Ferguson et al. (6) considered an ad hoc strategy to  
172 assist in overcoming this issue. The suggested approach involves the combination of this  
173 estimator with the conditional likelihood estimator, in a manner which is determined by the  
174 estimators' respective lengths of 95% confidence and credible intervals. Here, we instead  
175 investigated 5 alternative modifications to the original empirical Bayes method described  
176 above in order to better stabilize the tail of the estimated marginal density,  $\widehat{p(z)}$  and its  
177 derivative, particularly in the context of strong LD that is observed in high density  
178 genotyping arrays. These variations avoid the unappealing combination of two appreciably  
179 different estimators, empirical Bayes and conditional likelihood. The explored modifications  
180 were:

- 181 - Altering the minimum-BIC estimated spline function to be log-linear beyond the 10<sup>th</sup>  
182 largest negative and the 10<sup>th</sup> largest positive  $z$ -statistics

- 183 - Limiting the number of knots in the spline, in particular using 7 degrees of freedom as  
184 originally suggested by Efron (11)
- 185 - Utilizing smoothing splines, rather than natural splines, through the gam function in  
186 the R package mgcv (12), to avoid specifying knot positions, assuming a poisson  
187 distribution for the partition counts
- 188 - As above, but this time using a more realistic negative binomial distribution for these  
189 counts
- 190 - Employing splines with additional shape constraints, through the scam function in the  
191 R package scam (13), to enforce monotonicity of the estimated density function,  $\widehat{p}(z)$
- 192 More information about these modifications and their rationale is given in the supplementary  
193 material.

## 194 **FDR Inverse Quantile Transformation**

195 FDR Inverse Quantile Transformation (FIQT), as proposed by Bigdeli et al. (7),  
196 employs a straightforward two-step procedure in order to produce less biased association  
197 estimates. First, a FDR (false discovery rate) multiple testing adjustment is applied to the  $p$ -  
198 values of all SNPs, giving FDR adjusted  $p$ -values  $p_i^*$ ,  $i = 1, \dots, N$ . Following this, these  
199 adjusted  $p$ -values are transformed back to the  $z$ -statistic scale by means of an inverse  
200 Gaussian cumulative distribution function (cdf) and for each SNP, it is ensured that this new  
201  $z$ -statistic,  $\widehat{z}_i^*$ , has the same sign as its original effect size. Mathematically,  $\widehat{z}_i^*$ ,  $i = 1, \dots, N$ ,  
202 can be described as

$$\widehat{z}_i^* = \text{sign}(z_i) \Phi^{-1} \left( 1 - \frac{p_i^*}{2} \right). \quad (7)$$

203 For SNP  $i$ , its new estimated effect size is simply calculated as  $\widehat{\beta}_i = \widehat{z}_i^* \widehat{\text{se}}(\widehat{\beta}_i)$ .

204 The rationale that led to the use of this method is based on the analogy between  
205 performing multiple testing adjustments to  $p$ -values and reducing *Winner's Curse* bias in  
206 estimated SNP effect sizes, in which these effect sizes are in the form of  $z$ -statistics. In the  
207 attempt to correct for *Winner's Curse*, a shrinkage towards the null effect of zero is generally  
208 incurred by the  $z$ -statistics while the application of a multiple testing adjustment to  $p$ -values  
209 sees the growth of the  $p$ -values towards one, the null value.

210 This multiple testing adjustment is imposed through the implementation of the R  
211 function `p.adjust`. This is followed by the use of the R function `qnorm` for the purpose of  
212 back-transformation. However, near zero  $p$ -values can prove problematic when evaluating  
213 `qnorm` and thus, a restraint is incorporated in FIQT which results in the association estimates  
214 of SNPs with very large  $z$ -statistics, e.g. greater than 37, failing to be adjusted.

## 215 **Bootstrap**

216 Inspired by the bootstrap resampling method detailed in Sun et al. (14), we have  
217 established a similar approach which can be easily applied to published sets of GWAS  
218 summary statistics without requiring original individual-level data. In addition, a second  
219 advantage of our new method is a considerable improvement in computational efficiency  
220 over the method described in Sun et al. (14).

221 This procedure begins with arranging all  $N$  SNPs according to their original  $z$ -  
222 statistics,  $z_i = \frac{\hat{\beta}_i}{\text{se}(\hat{\beta}_i)}$ , in descending order, that is a labelling of SNPs is assumed such that  $z_1 >$   
223  $z_2 > \dots > z_N$ . A randomized estimate of the extent of ranking bias for the  $k^{\text{th}}$  largest  $z$ -statistic  
224 is calculated by means of the parametric bootstrap as follows:

- 225 1) A value  $\hat{\beta}_i^b$  is simulated for SNP  $i$ ,  $i = 1, \dots, N$ , independently, from a Gaussian  
226 distribution with mean  $\hat{\beta}_i$  and standard deviation  $\text{se}(\hat{\beta}_i)$ , i.e.

$$\hat{\beta}_i^b \sim N\left(\hat{\beta}_i, \widehat{\text{se}}(\hat{\beta}_i)\right). \quad (8)$$

227

228 2) Upon obtaining  $\hat{\beta}_i^b$  for  $i = 1, \dots, N$ , the  $z_i^b$ -statistic of SNP  $i$  is defined as

$$z_i^b = \frac{\hat{\beta}_i^b}{\widehat{\text{se}}(\hat{\beta}_i)}. \quad (9)$$

229 We define  $A(k)$  as the index corresponding to the  $k^{\text{th}}$  largest entry in the vector:

$$230 \quad [z_1^b, \dots, z_N^b] = \left[ \frac{\hat{\beta}_1^b}{\widehat{\text{se}}(\hat{\beta}_1)}, \dots, \frac{\hat{\beta}_N^b}{\widehat{\text{se}}(\hat{\beta}_N)} \right].$$

231 3) Then, the estimated bias of SNP  $k$ , the SNP with the  $k^{\text{th}}$  largest original  $z$ -statistic,

232 takes the following form:

$$\text{bias}_k = \frac{\hat{\beta}_{A(k)}^b - \hat{\beta}_{A(k)}^{\text{oob}}}{\widehat{\text{se}}(\hat{\beta}_{A(k)})} = \frac{\hat{\beta}_{A(k)}^b - \hat{\beta}_{A(k)}}{\widehat{\text{se}}(\hat{\beta}_{A(k)})}, \quad (10)$$

233 in which  $\hat{\beta}_{A(k)}^b$  is the bootstrap value of the SNP ranked in position  $k$  in the ordering of  
 234  $z_i^b$ -statistics,  $\hat{\beta}_{A(k)}^{\text{oob}} = \hat{\beta}_{A(k)}$  is that same SNP's original  $\beta$  estimate and  $\widehat{\text{se}}(\hat{\beta}_{A(k)})$  its  
 235 standard error.

236 In the next step of the process, a cubic smoothing spline is fitted to the data in which the  $z$ -  
 237 statistics are considered as the inputs and  $\text{bias}_k$ , their corresponding outputs. The predicted  
 238 values from this model fitting provides new estimates for the bias correction,  $\text{bias}_k^*$  for each  
 239 SNP. This additional stage in which  $\text{bias}_k^*$  is obtained reduces the need for more than one  
 240 bootstrap iteration for each SNP in order to ensure competitive performance of the method.  
 241 This results in a faster approach with increased accuracy. Finally, the new estimate for the  
 242 true effect size of SNP  $k$ , the SNP with the  $k^{\text{th}}$  largest original  $z$ -statistic, is defined as:  $\hat{\beta}_k^* =$

$$243 \quad \hat{\beta}_k - \widehat{\text{se}}(\hat{\beta}_k) \cdot \text{bias}_k^*.$$

244 In addition to those mentioned previously, there are several notable differences  
245 between our algorithm described above and the method proposed by Sun et al. (14). Firstly, it  
246 is the parametric bootstrap that is used here to estimate the magnitude of bias for each SNP as  
247 opposed to the more common nonparametric bootstrap which requires individual-level data.  
248 Our method draws only one bootstrap resample, i.e. only one bootstrap value  $\hat{\beta}_i^b$  is simulated  
249 for SNP  $i$ ,  $i = 1, \dots, N$ . It also includes an extra step which involves the use of a smoothing  
250 spline. In contrast, Sun et al. (14) express the need for a number of bootstrap samples, e.g. at  
251 least 100, in their approach. Furthermore, our algorithm based on the parametric bootstrap  
252 only corrects for ranking bias, and not threshold-selection bias.

## 253 **Simulation study**

254 The simulation study followed a factorial design in which GWAS summary statistics  
255 were simulated for a quantitative trait under 8 different genetic architectures, described by  
256 combinations of three parameters, namely sample size  $n$ , heritability  $h^2$ , polygenicity  
257 (proportion of effect SNPs)  $\pi$ . The following the values chosen for these parameters:

- 258 - sample size  $n \in \{30000, 300000\}$
- 259 - heritability  $h^2 \in \{0.3, 0.8\}$
- 260 - polygenicity  $\pi \in \{0.01, 0.001\}$

261 Assuming a selection coefficient equal to zero and a normal distribution of effect sizes, for a  
262 fixed array of  $N = 1,000,000$  SNPs, our strategy entailed imposing a simple correlation  
263 structure on the SNPs in order to imitate the presence of linkage disequilibrium (LD) in real  
264 data. It was assumed that the same correlation structure exists in independent blocks of 100  
265 SNPs. Thus, for each block of 100 SNPs, the estimated effect sizes,  $\hat{\beta}_i$  were simulated using:

$$\hat{\beta} \sim N\left(\mathbf{D}^{-\frac{1}{2}}\mathbf{R}\mathbf{D}^{\frac{1}{2}}\mathbf{b}, \mathbf{D}^{-\frac{1}{2}}\mathbf{R}\mathbf{D}^{-\frac{1}{2}}\sigma^2\right). \quad (11)$$

266 Here,  $\mathbf{b}$  is a vector containing the true SNP-trait effect sizes which have been scaled to ensure  
 267 that the phenotype has variance 1, i.e.  $\sigma^2 = 1$ . The matrix  $\mathbf{D}$  is a diagonal  $100 \times 100$  matrix, in  
 268 which  $d_i = n \cdot 2 \cdot \text{maf}_i(1 - \text{maf}_i)$  and  $\text{maf}_i$  is the minor allele frequency of SNP  $i$ , while  $\mathbf{R}$  is  
 269 a simple  $100 \times 100$  matrix of inter-genotype correlations, with  $R_{ij} = \hat{\rho}^{|i-j|}$  and  $\hat{\rho} = 0.9825$ .  
 270 The reasoning for the selection of this value for  $\hat{\rho}$  and why it was considered suitable, as well  
 271 as other details regarding this simulation, are described in the supplementary material. For  
 272 each SNP, values for  $\hat{\beta}_i$ ,  $\widehat{\text{se}}(\hat{\beta}_i)$  and  $E(\hat{\beta}_i)$  were produced with  $E(\hat{\beta}_i)$  obtained using  $E(\hat{\beta}) =$   
 273  $\mathbf{D}^{-\frac{1}{2}}\mathbf{R}\mathbf{D}^{\frac{1}{2}}\mathbf{b}$ . For each of these 8 different genetic architectures, 100 sets of summary statistics  
 274 were simulated.

275 The *Winner's Curse* correction methods detailed in 'Materials and methods' were  
 276 applied to each data set using the R package 'winnerscurse', producing adjusted estimated  
 277 effect sizes,  $\hat{\beta}_{\text{adj}, i}$ , for each SNP  $i$ ,  $i = 1, \dots, N$ . The performance of these methods were  
 278 investigated at two different significance thresholds, namely  $\alpha_1 = 5 \times 10^{-8}$  and  $\alpha_2 = 5 \times 10^{-4}$ ,  
 279 with a stronger focus given to the more commonly used genome-wide significance threshold  
 280 of  $\alpha_1 = 5 \times 10^{-8}$ . In order to assess each method's ability at providing less biased SNP-trait  
 281 association estimates, both the estimated change in mean squared error (MSE) and estimated  
 282 change in root mean squared error (RMSE) of significant SNPs due to method  
 283 implementation were computed for each data set and method. For simplicity, let  $i = 1, \dots, N_{\text{sig}}$   
 284 represent indexes for the significant SNPs in a particular simulated set of summary statistics,  
 285 i.e.  $N_{\text{sig}}$  is the number of SNPs which satisfy  $|z_i| > c$  with  $|z_i| = \left| \frac{\hat{\beta}_i}{\widehat{\text{se}}(\hat{\beta}_i)} \right|$ ,  $c = \Phi^{-1}\left(1 - \frac{\alpha}{2}\right)$  and  
 286  $\alpha \in \{\alpha_1, \alpha_2\}$ . Then, the estimated change in MSE of significant SNPs may be defined as:

$$\frac{1}{N_{\text{sig}}} \sum_{i=1}^{N_{\text{sig}}} (\hat{\beta}_{\text{adj},i} - \beta_i)^2 - \frac{1}{N_{\text{sig}}} \sum_{i=1}^{N_{\text{sig}}} (\hat{\beta}_i - \beta_i)^2 \quad (12)$$

287 while the estimated change in RMSE is defined similarly as:

$$\sqrt{\frac{1}{N_{\text{sig}}} \sum_{i=1}^{N_{\text{sig}}} (\hat{\beta}_{\text{adj},i} - \beta_i)^2} - \sqrt{\frac{1}{N_{\text{sig}}} \sum_{i=1}^{N_{\text{sig}}} (\hat{\beta}_i - \beta_i)^2}. \quad (13)$$

288 The change in MSE and RMSE for each method was calculated for only those data sets in  
289 which at least one significant SNP was detected. In addition to these two metrics, the relative  
290 change in MSE, which is equal to the change in MSE divided by the naïve MSE, was  
291 computed in a similar manner. For a given correction method, this value provides the  
292 percentage improvement in MSE due to applying that method to the set of summary statistics.

293 In addition to the above simulation set-up, GWAS summary statistics were simulated  
294 and methods evaluated under the assumption that SNPs were independent. In this instance,  
295 the study was extended to 24 genetic architectures in which the selection coefficient  $S$  took  
296 values -1 and 1 as well as 0. This simulation process which incorporates an independence  
297 assumption was repeated in a similar fashion for a binary trait with a normal distribution of  
298 effect sizes. Furthermore, a quantitative phenotype with a bimodal effect size distribution as  
299 well as one with a skewed distribution were also considered. In order to reduce computation  
300 time, only 50 sets of summary statistics were simulated for each combination of the four  
301 parameters for these three additional situations.

## 302 **Empirical analysis**

303 In order to compare the performance of these *Winner's Curse* correction methods with  
304 respect to real data, three different UK Biobank data sets were used, namely body mass index  
305 (BMI), height and type 2 diabetes (T2D). As with real data, the true effect size of each SNP is  
306 unknown and so it is more difficult to assess how much each method reduces the bias induced  
307 by *Winner's Curse*. To overcome this challenge, each original large data set was randomly

308 split in two, leaving between 166,172 and 166,687 individuals in each of the six smaller data  
309 sets. This provided the ability to execute two independent GWASs of similar sample size for  
310 each trait in which one GWAS was designated as the discovery GWAS and the other the  
311 replication GWAS. The unbiased replication GWAS association estimates can then be used  
312 as proxies for the true effect sizes of the SNPs found to be significant in the discovery  
313 GWAS. PLINK 2.0 (15) was used to perform quality control as well as each of the statistical  
314 analyses.

315 The required quality control steps which took place beforehand included the removal  
316 of related individuals. Samples which had been identified as outliers with respect to  
317 heterozygosity and missingness together with samples with discordant sex information and  
318 those suffering from chromosomal aneuploidy were also discarded. Furthermore, non-  
319 European samples which were identified by principal component analysis (PCA) using 1000  
320 Genomes data were removed. With respect to variants, only those with an information score  
321 greater than 0.8, a minor allele frequency greater than 0.01, a genotyping rate of at least 98%  
322 and those that passed the Hardy-Weinberg test at the specified significance threshold of  $1 \times$   
323  $10^{-8}$  were included.

324 The methods of interest were applied to the summary statistics of each discovery  
325 GWAS using the R package ‘winnerscurse’. Evaluation took place by computing the  
326 estimated MSE of  $N_{\text{sig}}$  significant SNPs in that GWAS, defined as:

$$\frac{1}{N_{\text{sig}}} \sum_{i=1}^{N_{\text{sig}}} (\hat{\beta}_{\text{disc,adj},i} - \hat{\beta}_{\text{rep},i})^2 - \frac{1}{N_{\text{sig}}} \sum_{i=1}^{N_{\text{sig}}} (\text{se}(\widehat{\beta}_{\text{disc},i}))^2. \quad (14)$$

327 For each of the three traits, it was possible to evaluate the performance of methods twice as in  
328 each case, the original roles of the two independent data sets, i.e. discovery and replication,  
329 could be switched and re-evaluation of methods could then take place with respect to the  
330 SNPs that were deemed significant in this new discovery GWAS.



## 331 **Results**

### 332 **Simulation study**

#### 333 **When is winner's curse bias most prominent?**

334 A simulation study in which a simple correlation structure was imposed on the set of  
335  $N = 1,000,000$  SNPs was first executed, as described in 'Materials and methods'. Before  
336 application of the *Winner's Curse* correction methods to the sets of summary statistics, an  
337 attempt to gain an insight into the simulation scenarios in which *Winner's Curse* bias is most  
338 prominent was made. This was done by computing the average number of significant SNPs,  
339 the average naïve MSE of significant SNPs and the average proportion of significant SNPs  
340 that had significantly overestimated effect sizes in each setting with respect to two  
341 significance thresholds,  $5 \times 10^{-8}$  and  $5 \times 10^{-4}$ . A SNP is defined as being significantly  
342 overestimated or as having a significantly more extreme effect size estimate if it satisfies the  
343 condition:

$$|\hat{\beta}_i| > |\beta_i| + 1.96 \cdot \widehat{\text{se}}(\hat{\beta}_i) \quad (15)$$

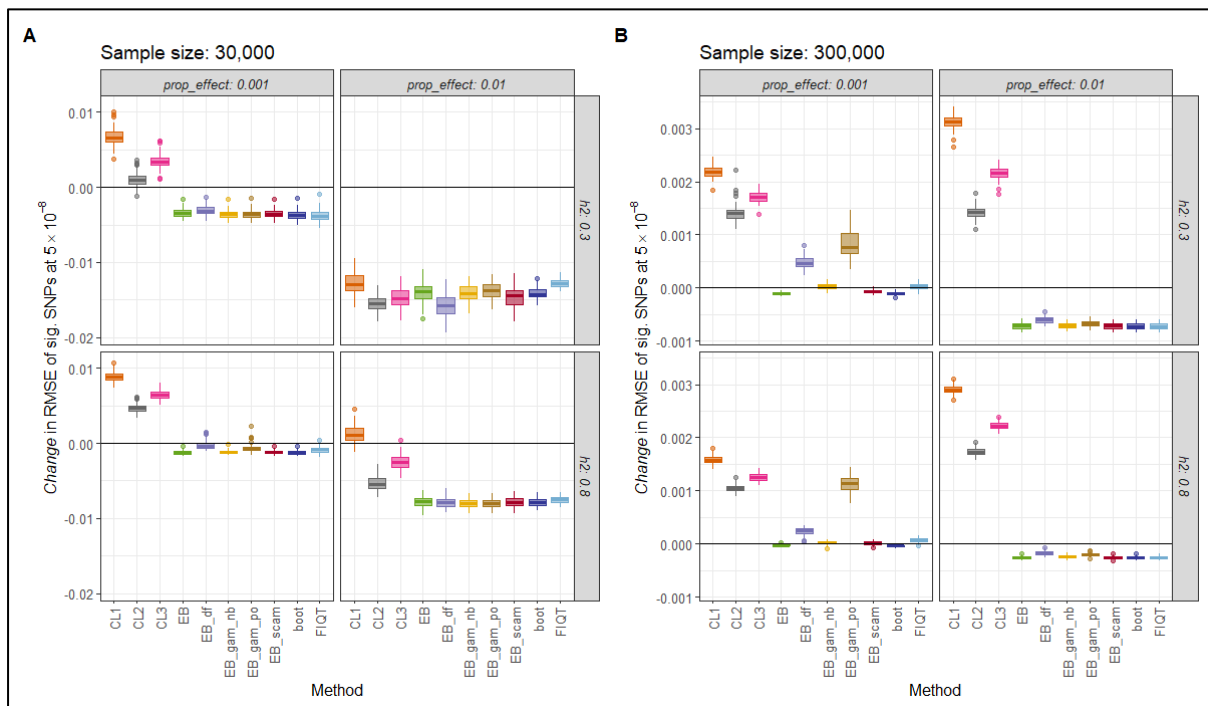
344 Thus, the proportion of significant SNPs that are significantly overestimated is considered to  
345 be representative of the proportion of significant SNPs with effect size estimates that greatly  
346 suffer from *Winner's Curse* bias. As detailed in S1 Table, it was clear that as sample size was  
347 increased from 30,000 to 300,000, this proportion of significant SNPs decreased. The other  
348 two parameters which played key roles in defining the various simulated genetic architectures  
349 were heritability and polygenicity. It was observed that the proportion of significantly  
350 overestimated significant SNPs decreased when heritability was increased from 0.3 to 0.8,  
351 but when the value representing trait polygenicity was increased from 0.001 to 0.01, this  
352 proportion decreased.

353 In fact, it was also noted that as the number of significant SNPs increased, both the  
354 MSE of significant SNPs and the proportion of these that were significantly overestimated  
355 decreased. This can be clearly seen in S1 Fig. This indicates that as the number of samples in  
356 a study and as the number of SNPs passing the significance threshold increases, bias induced  
357 by *Winner's Curse* will be less of an issue among significant SNPs. In terms of genetic  
358 architecture characteristics, these results suggest that the presence of *Winner's Curse* bias in  
359 the estimated effect sizes of significant SNPs should be of a greater concern when  
360 investigating traits with lower heritability or traits which have a larger proportion of effect  
361 SNPs.

### 362 **Evaluation of performance at $p < 5 \times 10^{-8}$**

363 With respect to the evaluation of methods, we focus on the results of computing the  
364 quantity 'change in RMSE over all significant SNPs due to method implementation' for each  
365 method, with obtaining a negative value being desirable. These results are provided in S2 and  
366 S5 Tables. At a threshold of  $5 \times 10^{-8}$ , several observations were notable. Firstly, for scenarios  
367 in which sample size has been designated the greater value of 300,000, the effect of applying  
368 the methods is on a much smaller scale to those scenarios with sample sizes of 30,000. This is  
369 evident from the large difference in the values on the y-axis between plots (A) and (B) of Fig  
370 1. This observation ties in with the fact that the magnitude of *Winner's Curse* bias is greater at  
371  $n = 30,000$ . At this  $5 \times 10^{-8}$  significance threshold, the conditional likelihood methods are  
372 seen to perform poorly, especially when sample sizes are increased to 300,000. In most  
373 instances, these methods provide worse association estimates than the naïve approach, often  
374 increasing the RMSE. The reason for this observation is over-correction of estimated effect  
375 sizes, especially those that lie close to the significance threshold.

376 **Fig 1. Estimated change in RMSE of significant SNPs at threshold  $5 \times 10^{-8}$  for each method and**  
 377 **simulation setting, with a simple correlation structure imposed on the set of SNPs.** The estimated  
 378 change in RMSE of significant SNPs at a threshold of  $5 \times 10^{-8}$  (y-axis), as defined by Eq (13), is  
 379 plotted for each correction method (x-axis), for each of the eight simulation settings. This figure  
 380 corresponds to simulation settings in which a simple correlation structure has been imposed on the set  
 381 of SNPs. Two four-panel plots, (A) and (B), are shown in the figure, in which (A) contains results  
 382 related to settings with a sample size of 30,000 and (B) contains results for sample sizes of 300,000.  
 383 The rows of these multi-panel plots represent heritability and the columns represent the proportion of  
 384 effect SNPs. Each panel contains a boxplot for each correction method. As 100 sets of summary  
 385 statistics were simulated for each simulation setting, an individual boxplot displays the distribution of  
 386 estimated change in RMSE values obtained across the 100 sets with respect to a particular method and  
 387 setting. The solid black line in each panel, representing no change in RMSE, is included in order to  
 388 highlight which methods consistently provide negative values for the estimated change in RMSE.  
 389 These methods are considered to be the best performing *Winner's Curse* correction methods.



390

391 The novel bootstrap method is one of the most consistent methods at providing less  
 392 biased SNP-trait association estimates at this threshold. In all situations depicted, it has one of

393 the largest negative values for the change in RMSE of significant SNPs and on average,  
394 improves the MSE by 26% across the 8 settings, as shown in S4 Table. FIQT tends to  
395 perform in a somewhat similar manner to this bootstrap method. With respect to the empirical  
396 Bayes method and its variations, for sample sizes of 300,000, the best performing versions  
397 were the original empirical Bayes method, ‘EB-gam-nb’ and ‘EB-scsm’. Note that the  
398 original empirical Bayes method is the form most similar to that proposed in Ferguson et al.  
399 (6) which also includes a restriction on the tails of the distribution of  $z$ -statistics. With the  
400 lower sample size, all variations perform very similarly. However, the empirical Bayes  
401 variants ‘EB-df’ and ‘EB-gam-po’ performed less well overall. In the case of ‘EB-gam-po’,  
402 this might be partly due to convergence problems that sometimes occurred in obtaining the  
403 poisson regression fit. In general, we advise caution in utilizing the empirical Bayes results in  
404 the context of convergence warnings from R.

#### 405 **Evaluation of performance at $p < 5 \times 10^{-4}$**

406 A lower threshold of  $5 \times 10^{-4}$  was also investigated. Less emphasis is placed on the  
407 results obtained at this threshold as it is possible that many false positives are detected here,  
408 i.e. many SNPs that in fact have a true effect size of zero pass the significance threshold,  
409 although lower thresholds may be useful for the construction of polygenic risk scores.  
410 Therefore, as these *Winner’s Curse* correction methods are all considered to be shrinkage  
411 methods, improvements in the RMSE over all significant SNPs would be expected. However,  
412 as can be seen in S3 Fig, positive values are witnessed at this threshold when sample sizes are  
413 large. However, these positive values most often occur for the conditional likelihood methods  
414 which seem to be the worst performers overall. It seems here that the most consistent and best  
415 performing methods are the bootstrap method, the original empirical Bayes method and the  
416 empirical Bayes method which uses shape constrained additive models (SCAMs). These

417 methods all reduce the MSE of significant SNPs at this threshold by an average of at least  
418 32%, as shown in S7 Table.

### 419 **Additional simulations in absence of Linkage Disequilibrium**

420 In order to demonstrate potential performance of the *Winner's Curse* methods in the  
421 context of SNP-trait associations from genome-wide arrays with lower SNP density or LD-  
422 pruned datasets, we also examined the less complex situation in which SNPs are independent.  
423 The results of these extra simulations are shown in S4-S11 Figs and described in depth in the  
424 supplementary material. In this setting, with a normal effect size distribution, the most  
425 consistent methods in terms of reducing the RMSE of significant SNPs were the original  
426 empirical Bayes method and 'EB-gam-nb', the variation of the empirical Bayes method  
427 which employs smoothing splines and assumes a negative binomial count distribution. Just as  
428 was observed in the simulations with linkage disequilibrium, the conditional likelihood  
429 methods perform poorly and often result in an increase in the evaluation metric in comparison  
430 to the naïve approach, while the proposed bootstrap method continued to exhibit competitive  
431 performance.

### 432 **Empirical analysis**

433 The results of an initial exploration of the six UK Biobank sets of summary statistics  
434 are detailed in Table 1. From trait to trait, there is a large difference in the number of SNPs  
435 with  $p$ -values lower than the genome-wide significance threshold of  $5 \times 10^{-8}$ . Values for the  
436 proportion of these SNPs with significantly overestimated effect sizes in each discovery  
437 GWAS are included. A comparison of BMI and height GWASs at the  $5 \times 10^{-8}$  threshold tends  
438 to indicate that as the number of significant SNPs increases, the proportion that are  
439 significantly overestimated decreases. This trend is even more apparent at the larger threshold  
440 of  $5 \times 10^{-4}$ , and as stated above, was also clearly observed in the simulated data.

441 **Table 1. The number of significant SNPs at two significance thresholds,  $5 \times 10^{-8}$  and  $5 \times 10^{-4}$ ,**  
 442 **with proportions that indicate the extent of *Winner's Curse* bias for each data set.**

GWAS	No. sig. SNPs ( $5 \times 10^{-8}$ )	Prop. sig. SNPs with <i>smaller</i> replication estimate ( $5 \times 10^{-8}$ )	Prop. sig. SNPs <i>significantly</i> overestimated ( $5 \times 10^{-8}$ )	No. sig. SNPs ( $5 \times 10^{-4}$ )	Prop. sig. SNPs with <i>smaller</i> replication estimate ( $5 \times 10^{-4}$ )	Prop. sig. SNPs <i>significantly</i> overestimated ( $5 \times 10^{-4}$ )
<b>BMI 1</b>	6,908	0.7135	0.2251	94,173	0.8365	0.3386
<b>BMI 2</b>	7,951	0.8009	0.3202	98,351	0.8455	0.3604
<b>T2D 1</b>	31	0.0645 <sup>a</sup>	0.0645	5,832	0.9830	0.8433
<b>T2D 2</b>	76	1	0.1579	5,507	0.9951	0.8397
<b>Height 1</b>	70,020	0.6444	0.1829	257,000	0.6940	0.2095
<b>Height 2</b>	70,634	0.6824	0.1772	268,497	0.7179	0.2406

443 Table 1 details an initial exploration of the six UK Biobank data sets. The number of significant SNPs  
 444 identified for each data set at two significance thresholds,  $5 \times 10^{-8}$  and  $5 \times 10^{-4}$ , is provided in the  
 445 table. The table also contains the proportion of these significant SNPs that were seen to have smaller  
 446 estimated effect sizes, in terms of absolute value, in their respective replication GWAS. The final  
 447 column for each threshold provides the proportion of significant SNPs that have significantly  
 448 overestimated effect sizes. The estimated effect size of a SNP has been defined as significantly  
 449 overestimated according to Eq (15), but in which the true effect size is replaced by the estimated  
 450 effect size obtained in the corresponding replication GWAS. These values give an indication of the  
 451 extent of *Winner's Curse* bias present for each data set and threshold.

452 <sup>a</sup>When the first T2D data set was used for the discovery GWAS, most of the 31 significant SNPs had  
 453 larger estimated effect sizes in the replication data set. Discussion of this atypical observation is  
 454 included in the main text.

455

## 456 **The problem of Linkage Disequilibrium in real data**

457 Naturally, the results of engagement with real data sets are more complex than those  
 458 of the simulation study. For example, it was noted that for BMI, in one instance, all  
 459 significant SNPs which had a  $z$ -value greater than 15 in the discovery GWAS had association

460 estimates in the replication GWAS which were in fact greater. This observation can be  
461 clearly seen in S12 Fig in which  $z$ -statistics are plotted against estimated bias for each data  
462 set, with estimated bias of SNP  $i$  defined as:

$$\widehat{\text{bias}}_i = \hat{\beta}_{\text{disc},i} - \hat{\beta}_{\text{rep},i} \quad (16)$$

463 This finding is of course contrary to what is expected. However, these SNPs with  $z$ -values  
464 greater than 15 were all in strong linkage disequilibrium and thus, represented a single  
465 independent signal. It can be seen in Table 1 that a similar result was noted when the first  
466 T2D data set was used as the discovery GWAS. When using a significance threshold of  $5 \times$   
467  $10^{-8}$ , most of the 31 significant SNPs had larger estimated effect sizes in the replication  
468 GWAS than in the discovery GWAS. In these cases, we need to be careful not to over-  
469 generalize or interpret the results of applying a *Winner's Curse* correction, given that there  
470 may be very few independent association signals at  $p < 5 \times 10^{-8}$ .

### 471 **Evaluation of performance at $p < 5 \times 10^{-8}$**

472 As stated in 'Materials and methods', the methods were evaluated using the estimated  
473 MSE of SNPs which passed the chosen significance threshold. Using the threshold of  $5 \times 10^{-8}$ ,  
474 the estimated MSE for each method and GWAS combination are displayed in Table 2  
475 while Fig 2 provides a corresponding illustration of these values. In this figure, the light blue  
476 bar as well as the black dotted horizontal line mark the estimated MSE obtained using the  
477 naïve approach, i.e. when no *Winner's Curse* correction method has been applied and the raw  
478 effect estimates are used. This provides a standard to which the performance of each method  
479 can be directly compared with, in which it is desired that method application will result in an  
480 estimated MSE less than this approach. Similar to the section above describing the results of  
481 the simulation study, the poor performance of the conditional likelihood methods is evident.

482 In 5 out of the 6 independent instances, it was observed that at least one of these methods had  
 483 a greater estimated MSE than that of the naïve approach.

484 **Table 2. Estimated MSE of significant SNPs at threshold  $5 \times 10^{-8}$  for each method and data set.**

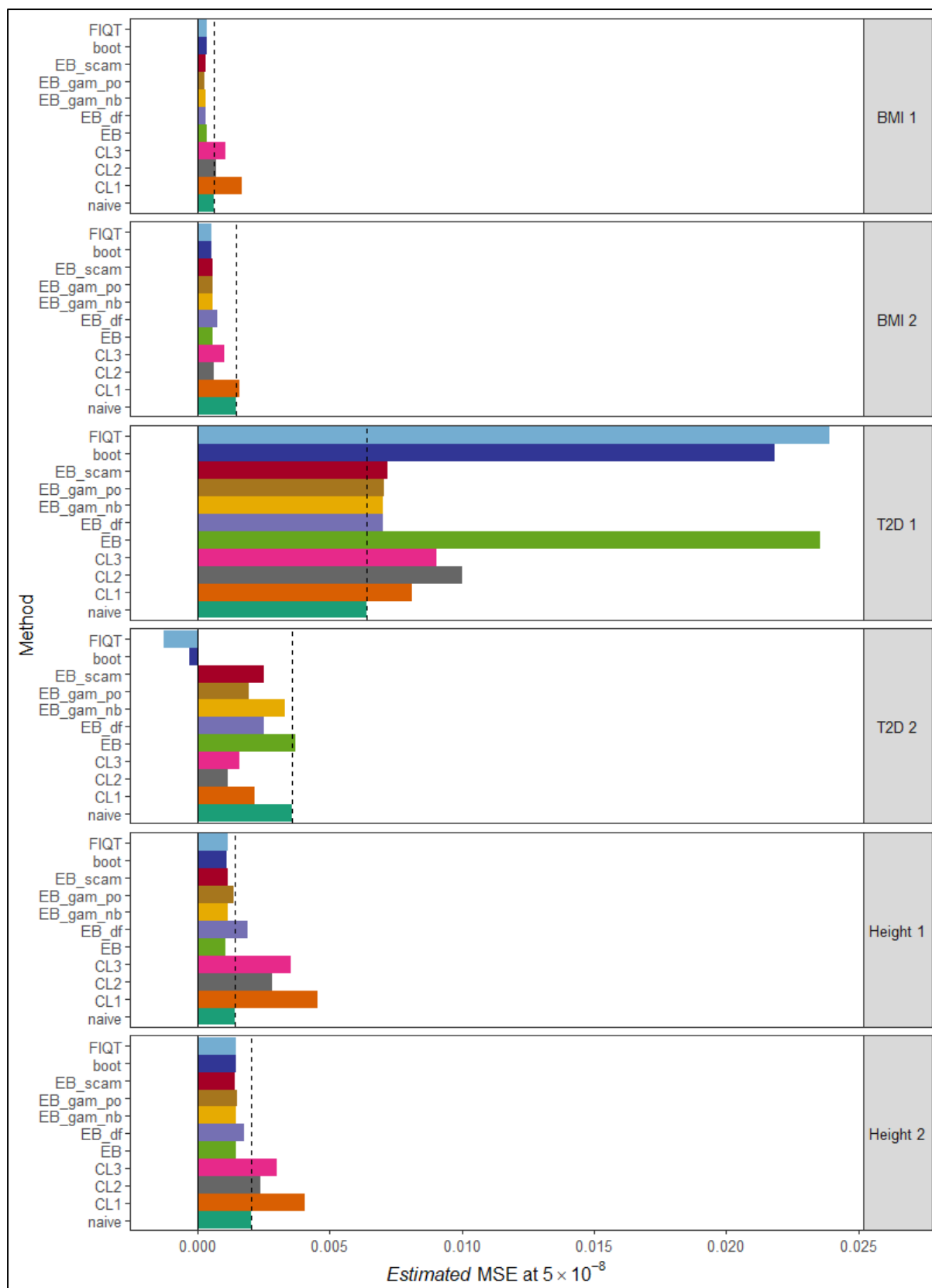
GWAS	BMI 1	BMI 2	T2D 1	T2D 2	Height 1	Height 2
naive	0.00065	0.00148	0.00641	0.00358	0.00142	0.00206
CL1	0.00167	0.00162	0.00811	0.00217	0.00457	0.00408
CL2	0.00071	0.00065	0.01004	0.00116	0.00285	0.00239
CL3	0.00108	0.00103	0.00903	0.00159	0.00355	0.00303
EB	0.00036	0.00058	0.02354	0.00371	0.00109	0.00146
EB df=7	0.00031	0.00077	0.00701	0.00251	0.00189	0.00179
EB scam	0.00033	0.00058	0.00719	0.00252	0.00116	0.00141
EB gam-po	0.00029	0.00061	0.00706	0.00194	0.0014	0.00153
EB-gam-nb	0.00031	0.00059	0.00703	0.0033	0.00118	0.00145
boot	0.00034	0.00057	0.02103	-0.0003	0.00111	0.00148
FIQT	0.00039	0.00056	0.0239	-0.0013	0.00115	0.00149

485 Table 2 provides values for the estimated MSE of significant SNPs, as defined by Eq (14), using a  
 486 threshold of  $5 \times 10^{-8}$ , for each *Winner's Curse* correction method and UK Biobank data set. The first  
 487 row of values represents the estimated MSE obtained if the unadjusted estimated effect sizes of the  
 488 discovery GWAS are used and no correction method has been applied. This is followed by rows  
 489 which are representative of the use of different correction methods, i.e. the conditional likelihood  
 490 based methods, the empirical Bayes method and its variations, the proposed bootstrap method and  
 491 FIQT, respectively. As it is desirable to obtain lower estimated MSE values upon application of a  
 492 method, values which are greater than their corresponding naïve value have been shaded in grey. The  
 493 light green shaded cells highlight the method which resulted in the lowest estimated MSE value for  
 494 each data set.

495 **Fig 2. Estimated MSE of significant SNPs at threshold  $5 \times 10^{-8}$  for each method and data set.**

496 The estimated MSE of significant SNPs at a threshold of  $5 \times 10^{-8}$  (x-axis), as defined by Eq (14), is  
 497 plotted for each correction method (y-axis), for each of the six data sets. The estimated MSE obtained  
 498 for the naïve approach, when no *Winner's Curse* correction method is applied, is included,  
 499 represented by the darker green bar. The dashed black line also represents this value, in order to  
 500 highlight which methods provide estimated MSE values greater or less than that of the naïve  
 501 approach. All estimated MSE values plotted here are provided in Table 2.





502

503 On real data, the original empirical Bayes method proposed by Ferguson et al. (6)

504 performed poorly, sometimes failing to adjust estimated associations downwards. This

505 observation motivated the proposal of possible modifications as mentioned in ‘Materials and  
506 methods’. These suggested improvements, in particular the inclusion of the shape-constrained  
507 additive models and the use of the generalized additive model, resulted in slightly more  
508 consistent reductions in MSE over significant SNPs. In fact, taking all six data sets into  
509 account, it is ‘EB-scsm’ and ‘EB-gam-po’, which tend to be the best performing methods,  
510 having an average improvement on estimated MSE of greater than 29.4% over the naïve  
511 approach.

512         However, as stated previously, we must be cautious when using the first T2D data set  
513 to evaluate methods, and also when using the second T2D data set. The problem with linkage  
514 disequilibrium and very few independent signals is common to both data sets. In Fig 2 for the  
515 first T2D data set, it is witnessed that all methods result in greater estimated MSE values than  
516 the naïve approach, with the original empirical Bayes method, bootstrap and FIQT clearly  
517 greatly shrinking the estimated effect sizes of significant SNPs away from those larger  
518 replication effect sizes. Therefore, if we exclude these two T2D data sets and re-compute the  
519 average improvement in estimated MSE for each method, it is our proposed bootstrap method  
520 which is seen to be the dominant method with an average improvement of approximately  
521 40.2%.

## 522 **Evaluation of performance at $p < 5 \times 10^{-4}$**

523         This evaluation procedure was repeated using a larger significance threshold of  $5 \times$   
524  $10^{-4}$ . The results of which can be found summarised in S8 Table and Fig 3. At this threshold,  
525 for all 6 data sets, all of the methods produce estimated MSE values less than the naïve  
526 approach. Each version of the empirical Bayes method along with the bootstrap and FIQT  
527 lead to an average improvement in estimated MSE of between 65 and 70% with the

528 implementation of the empirical Bayes algorithm which incorporates shape constrained  
529 additive models (SCAMs) having the greatest average improvement of just over 70%.

530 **Fig 3. Estimated MSE of significant SNPs at threshold  $5 \times 10^{-4}$  for each method and data set.** The

531 estimated MSE of significant SNPs at a threshold of  $5 \times 10^{-4}$  (x-axis), as defined by Eq (14), is

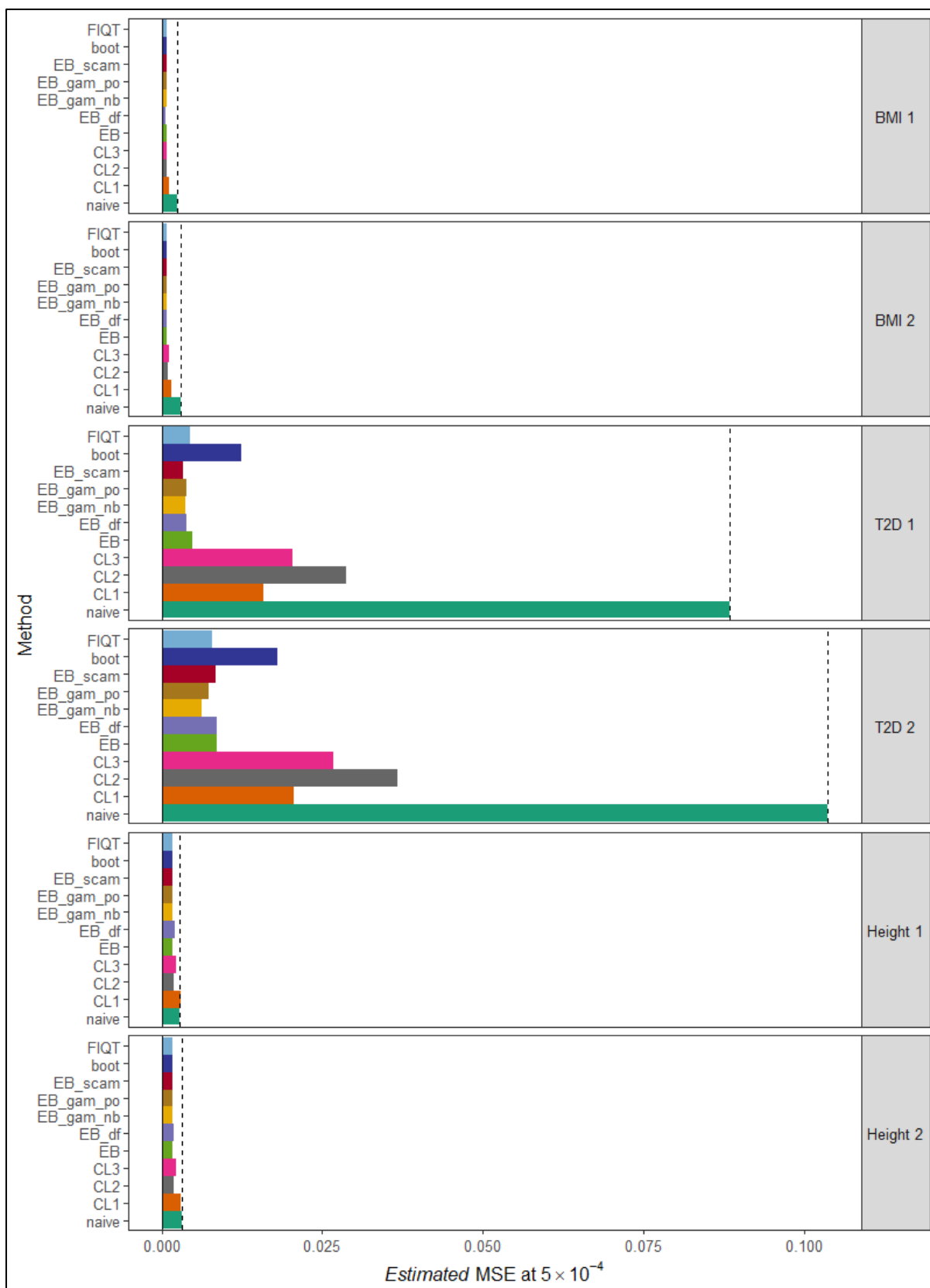
532 plotted for each correction method (y-axis), for each of the six data sets. The estimated MSE obtained

533 for the naïve approach, when no *Winner's Curse* correction method is applied, is included,

534 represented by the darker green bar. The dashed black line also represents this value, in order to

535 highlight which methods provide estimated MSE values greater or less than that of the naïve

536 approach. The estimated MSE values plotted here are provided in S6 Table.



537

## 538 Discussion

539           In this article, we investigated the problem of *Winner's Curse* bias which results in  
540 the estimated effect sizes of significant SNPs often being greater than their true values. Our  
541 work concentrated on methods that could be used to reduce this bias in settings in which only  
542 summary statistics of the GWAS that discovered these SNP-trait associations were available.  
543 We chose to focus on this particular situation as *Winner's Curse* correction methods which  
544 only require GWAS summary data tend to be very computational efficient and furthermore,  
545 this summary data is often much easier to access than individual-level data.

546           We performed a thorough evaluation and comparison of these methods using both  
547 simulated and real data sets. Our simulation study considered a wide range of genetic  
548 architectures including data sets in which a simple correlation structure had been imposed on  
549 the set of SNPs as well as data sets of independent SNPs. In addition, three UKBB real data  
550 sets were used for method evaluation purposes. As well as assessing currently published  
551 correction methods, we also explored several possible modifications that could be made in  
552 order to improve these methods. In particular, we looked at a number of variations of the  
553 empirical Bayes method and proposed an additional approach which uses the parametric  
554 bootstrap in order to establish suitable corrections for the estimated effect size of each SNP.  
555 The estimated mean squared error (MSE) was chosen as an appropriate metric in order to  
556 compare the methods. Due to the notable lack of software for implementation of *Winner's*  
557 *Curse* correction methods, we developed an R package, 'winnerscurse', as an accompaniment  
558 to the work described in this paper. This allows users to apply all methods discussed here, as  
559 well as the proposed modifications, to their sets of GWAS summary data.

560           As a first step in both our simulation study and engagement with real data, we  
561 computed the proportion of significant SNPs that were significantly overestimated and  
562 observed the common trend that as the number of SNPs passing the significance threshold  
563 increased, the proportion of those that were significantly overestimated decreased. This aligns

564 with the postulation that as sample sizes increase, *Winner's Curse* bias becomes less of a  
565 concern although it still exists. However, caution must be taken when working with real data  
566 sets, especially those of binary traits, in which a very small number of SNPs have been  
567 deemed significant at a certain threshold. In this instance, it may be that the significant SNPs  
568 are representative of only one or two independent signals. For example, in our first T2D data  
569 set, while using a threshold of  $5 \times 10^{-8}$ , we witnessed 93.5% of significant SNPs having  
570 greater replication effect size estimates than those obtained in the discovery GWAS.  
571 Fortunately, as sample sizes increase in the future and different diseases have greater  
572 numbers of true signals captured by their respective sets of significant SNPs, this issue will  
573 only be seen to present itself in rare circumstances.

574         With respect to method performance, it was clear that the conditional likelihood  
575 methods performed poorly as in most instances, especially for  $p < 5 \times 10^{-8}$ , these methods  
576 resulted in greater values for the estimated MSE among significant SNPs than the naïve  
577 approach. The other considered methods behaved much more similarly to the extent that we  
578 cannot state that there is a clear advantage of one method over another. Thus, the choice of  
579 which method a user should apply to their set of GWAS summary statistics in order to correct  
580 for *Winner's Curse* is dependent on personal preference. However, it is advised that when  
581 doing so, the possible limitations of the chosen method are understood well. Notably, the  
582 empirical Bayes methods have a clear theoretical advantage, but their performance can be  
583 restricted due to inaccurate estimation of the extreme tails of the  $z$ -statistic distribution. This  
584 estimation difficulty is particularly problematic when the existence of strong linkage  
585 disequilibrium results in clusters of associations in the tails. These clusters can be falsely  
586 detected as local modes in the distribution by automatic fitting algorithms. Some progress on  
587 improving estimation in the tails has been made here with the proposal of modifications that  
588 employ generalized additive models or shape constrained additive models. However, these

589 adaptations have not resulted in large enough improvements in order to claim objective  
590 superiority of the empirical Bayes methods over other approaches such as the bootstrap  
591 method or FIQT. In a setting in which the distribution of effect sizes is asymmetric, methods  
592 like the empirical Bayes and bootstrap, where the correction rule is not a function of absolute  
593 value z-statistics, possess the potential to perform better than FIQT and conditional  
594 likelihood methods. In spite of this fact, no tangible evidence of improved performance over  
595 FIQT was observed on the real data sets that we examined.

596         With both our set of simulations and real data analysis, we have aimed to be as  
597 comprehensive as possible as it is possible that differing method performance results may  
598 occur under differing genetic architectures, but this is an obviously difficult task. Informing  
599 these simulations appropriately is particularly challenging, especially when attempting to  
600 define the true effect size distribution. However, under the assumption of independent SNPs,  
601 we also investigated scenarios which had a bimodal or skewed distribution of effect sizes, as  
602 described in the supplementary material. Furthermore, for simulations involving correlated  
603 SNPs, we have assumed a very simplistic linkage disequilibrium (LD) structure in which the  
604 minor allele frequencies have been simulated independently of this LD structure. In contrast,  
605 the use of real data permitted the analysis of method performance in a realistic setting where  
606 a large degree of LD exists. However, this was limited to only three UKBB data sets. In the  
607 case of the binary trait T2D, it must be noted that due to the very small number of significant  
608 SNPs at  $p < 5 \times 10^{-8}$ , the results are deemed rather questionable here.

609         Due to space considerations, *Winner's Curse* correction methods which require both a  
610 discovery and replication GWAS in order to make suitable adjustments to estimated effect  
611 sizes have not been examined in this manuscript, even though several of these methods have  
612 been implemented in our developed R package, 'winnerscurse'. Furthermore, computation of  
613 standard errors of the adjusted estimated effect sizes have not been considered here.

614 However, for methods such as the empirical Bayes, bootstrap and FIQT, the R package,  
615 ‘winnerscurse’, utilizes the parametric bootstrap in order to obtain these standard errors. This  
616 package can also be used to provide confidence intervals for estimated effect sizes which  
617 have been corrected for *Winner’s Curse* using the conditional likelihood methods. In two-  
618 sample Mendelian randomization, it is known that as this *Winner’s Curse* bias can be present  
619 in the estimated SNP-exposure associations, the causal estimate will then suffer from bias.  
620 Thus, the *Winner’s Curse* correction methods explored in this paper can also be potentially  
621 used as plug-in corrections for two-sample MR. In addition, these methods could prove  
622 beneficial in the computation of polygenic risk scores, in order to reduce the effect of  
623 *Winner’s Curse* bias there.

## 624 **References**

- 625 1. Dudbridge F, Newcombe P. Replication and Meta-analysis of Genome-Wide  
626 Association Studies. *Handbook of Statistical Genomics: Two Volume Set*. 2019:631-50.
- 627 2. Jiang T, Gill D, Butterworth AS, Burgess S. An empirical investigation into the  
628 impact of winner’s curse on estimates from Mendelian randomization. medRxiv:  
629 2022.08.05.22278470 [Preprint]. 2022 [cited 2022 Oct 25]
- 630 3. Sadreev II, Elsworth BL, Mitchell RE, Peternoster L, Sanderson E, Davies NM, et al.  
631 Navigating sample overlap, winner’s curse and weak instrument bias in Mendelian  
632 randomization studies using the UK Biobank. medRxiv: 2021.06.28.21259622 [Preprint].  
633 [cited 2022 Oct 25]
- 634 4. Ruan Y, Choi SW, O’Reilly P. Investigating shrinkage methods to improve accuracy  
635 of GWAS and PRS effect size estimates. *European Neuropsychopharmacology*.  
636 2019;29(3):896-7.



- 637 5. Ghosh A, Zou F, Wright FA. Estimating odds ratios in genome scans: an approximate  
638 conditional likelihood approach. *The American Journal of Human Genetics*.  
639 2008;82(5):1064-74.
- 640 6. Ferguson JP, Cho JH, Yang C, Zhao H. Empirical Bayes correction for the Winner's  
641 Curse in genetic association studies. *Genetic epidemiology*. 2013;37(1):60-8.
- 642 7. Bigdeli TB, Lee D, Webb BT, Riley BP, Vladimirov VI, Fanous AH, et al. A simple  
643 yet accurate correction for winner's curse can predict signals discovered in much larger  
644 genome scans. *Bioinformatics*. 2016;32(17):2598-603.
- 645 8. Faye LL, Sun L, Dimitromanolakis A, Bull SB. A flexible genome-wide bootstrap  
646 method that accounts for ranking and threshold-selection bias in GWAS interpretation and  
647 replication study design. *Statistics in medicine*. 2011;30(15):1898-912.
- 648 9. Bowden J, Dudbridge F. Unbiased estimation of odds ratios: combining genomewide  
649 association scans with replication studies. *Genetic Epidemiology: The Official Publication of*  
650 *the International Genetic Epidemiology Society*. 2009;33(5):406-18.
- 651 10. Zhong H, Prentice RL. Bias-reduced estimators and confidence intervals for odds  
652 ratios in genome-wide association studies. *Biostatistics*. 2008;9(4):621-34.
- 653 11. Efron B. Tweedie's formula and selection bias. *Journal of the American Statistical*  
654 *Association*. 2011;106(496):1602-14.
- 655 12. Wood SN. Fast stable restricted maximum likelihood and marginal likelihood  
656 estimation of semiparametric generalized linear models. *Journal of the Royal Statistical*  
657 *Society: Series B (Statistical Methodology)*. 2011;73(1):3-36.
- 658 13. Pya N, Wood SN. Shape constrained additive models. *Statistics and computing*.  
659 2015;25(3):543-59.

- 660 14. Sun L, Dimitromanolakis A, Faye LL, Paterson AD, Waggott D, Bull SB. BR-  
661 squared: a practical solution to the winner's curse in genome-wide scans. *Human genetics*.  
662 2011;129(5):545-52.
- 663 15. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D, et al. PLINK:  
664 a tool set for whole-genome association and population-based linkage analyses. *The*  
665 *American journal of human genetics*. 2007;81(3):559-75.
- 666 16. Bosch E, Laayouni H, Morcillo-Suarez C, Casals F, Moreno-Estrada A, Ferrer-  
667 Admetlla A, et al. Decay of linkage disequilibrium within genes across HGDP-CEPH human  
668 samples: most population isolates do not show increased LD. *BMC genomics*. 2009;10(1):1-  
669 9.

## 670 **Supporting information**

671 **S1 Fig. Number of significant SNPs at threshold  $5 \times 10^{-8}$  plotted against the proportion of those**  
672 **SNPs with significantly overestimated effect sizes, for each simulation setting with a simple**  
673 **correlation structure imposed on the set of SNPs.** For 100 iterations of each of the eight simulation  
674 settings, the number of significant SNPs using a significance threshold of  $5 \times 10^{-8}$  (x-axis) is plotted  
675 against the number of these SNPs with significantly overestimated effect sizes (y-axis). The estimated  
676 effect size of a SNP has been defined as significantly overestimated according to Eq (15) in the main  
677 text. This figure corresponds to simulation settings in which a simple correlation structure has been  
678 imposed on the set of SNPs. These 8 different simulated genetic architectures are defined by  
679 combinations of three parameters, sample size  $n$ , heritability  $h^2$  and polygenicity  $\pi$ . The parameter  
680 values that have been chosen for each simulated scenario are shown in the table, while the legend at  
681 the bottom of the plot indicates which colour corresponds to which scenario.

682 **S2 Fig. Z-statistics plotted against bias for each simulation setting, with a simple correlation**  
683 **structure imposed on the set of SNPs.** For a single example of each of the eight simulation settings,

684  $z$ -statistics (x-axis) are plotted against bias (y-axis). The  $z$ -statistic of a SNP is defined as its estimated  
685 effect size divided by the standard error of that estimated effect size while the bias of a SNP is defined  
686 as its true effect size subtracted from its estimated effect size. This figure corresponds to simulation  
687 settings in which a simple correlation structure has been imposed on the set of SNPs. These 8  
688 different simulated genetic architectures are defined by combinations of three parameters, sample size  
689  $n$ , heritability  $h^2$  and polygenicity  $\pi$ . The parameter values that have been chosen for each simulated  
690 setting are shown in the subtitle of each plot. In each plot, the dark red dashed vertical line represents  
691 the  $z$ -statistic corresponding to a  $p$ -value of  $5 \times 10^{-8}$  and thus, any points outside these two dark red  
692 lines are SNPs with  $p$ -values passing the genome-wide significance threshold of  $5 \times 10^{-8}$ . In a similar  
693 manner, the light red dashed vertical line represents the greater significance threshold of  $5 \times 10^{-4}$ . The  
694 dark grey points highlight SNPs that have  $p$ -values less than  $5 \times 10^{-4}$  and have significantly  
695 overestimated effect sizes, as defined by Eq (15) in the main text.

696 **S3 Fig. Estimated change in RMSE of significant SNPs at threshold  $5 \times 10^{-4}$  for each method and**  
697 **simulation setting, with a simple correlation structure imposed on the set of SNPs.** The estimated  
698 change in RMSE of significant SNPs at a threshold of  $5 \times 10^{-4}$  (y-axis), as defined by Eq (13) in the  
699 main text, is plotted for each correction method (x-axis), for each of the eight simulation settings. This  
700 figure corresponds to simulation settings in which a simple correlation structure has been imposed on  
701 the set of SNPs.

702 **S4 Fig. Estimated change in RMSE of significant SNPs at threshold  $5 \times 10^{-8}$  for each method and**  
703 **simulation setting, assuming a quantitative trait, independent SNPs and a normal effect size**  
704 **distribution.** The estimated change in RMSE of significant SNPs at a threshold of  $5 \times 10^{-8}$  (y-axis),  
705 as defined by Eq (13) in the main text, is plotted for each correction method (x-axis), for each of the  
706 eight simulation settings with a selection coefficient of zero. This figure corresponds to simulation  
707 settings in which it is assumed that the trait of interest is quantitative, SNPs are independent and the  
708 effect sizes follow a normal distribution.

709 **S5 Fig. Estimated change in RMSE of significant SNPs at threshold  $5 \times 10^{-4}$  for each method and**  
710 **simulation setting, assuming a quantitative trait, independent SNPs and a normal effect size**  
711 **distribution.** The estimated change in RMSE of significant SNPs at a threshold of  $5 \times 10^{-4}$  (y-axis),  
712 as defined by Eq (13) in the main text, is plotted for each correction method (x-axis), for each of the  
713 eight simulation settings with a selection coefficient of zero. This figure corresponds to simulation  
714 settings in which it is assumed that the trait of interest is quantitative, SNPs are independent and the  
715 effect sizes follow a normal distribution.

716 **S6 Fig. Estimated change in RMSE of significant SNPs at threshold  $5 \times 10^{-8}$  for each method and**  
717 **simulation setting, assuming a quantitative trait, independent SNPs and a bimodal effect size**  
718 **distribution.** The estimated change in RMSE of significant SNPs at a threshold of  $5 \times 10^{-8}$  (y-axis),  
719 as defined by Eq (13) in the main text, is plotted for each correction method (x-axis), for each of the  
720 eight simulation settings with a selection coefficient of zero. This figure corresponds to simulation  
721 settings in which it is assumed that the trait of interest is quantitative, SNPs are independent and the  
722 effect sizes follow a bimodal distribution. Note that for this simulation set-up, the alternative  
723 variations of the empirical Bayes method have been excluded and only 50 sets of summary statistics  
724 were simulated for each setting.

725 **S7 Fig. Estimated change in RMSE of significant SNPs at threshold  $5 \times 10^{-4}$  for each method and**  
726 **simulation setting, assuming a quantitative trait, independent SNPs and a bimodal effect size**  
727 **distribution.** The estimated change in RMSE of significant SNPs at a threshold of  $5 \times 10^{-4}$  (y-axis),  
728 as defined by Eq (13) in the main text, is plotted for each correction method (x-axis), for each of the  
729 eight simulation settings with a selection coefficient of zero. This figure corresponds to simulation  
730 settings in which it is assumed that the trait of interest is quantitative, SNPs are independent and the  
731 effect sizes follow a bimodal distribution. Note that for this simulation set-up, the alternative  
732 variations of the empirical Bayes method have been excluded and only 50 sets of summary statistics  
733 were simulated for each setting.

734 **S8 Fig. Estimated change in RMSE of significant SNPs at threshold  $5 \times 10^{-8}$  for each method and**  
735 **simulation setting, assuming a quantitative trait, independent SNPs and a skewed effect size**  
736 **distribution.** The estimated change in RMSE of significant SNPs at a threshold of  $5 \times 10^{-8}$  (y-axis),  
737 as defined by Eq (13) in the main text, is plotted for each correction method (x-axis), for each of the  
738 eight simulation settings with a selection coefficient of zero. This figure corresponds to simulation  
739 settings in which it is assumed that the trait of interest is quantitative, SNPs are independent and the  
740 effect sizes follow a skewed distribution. Note that for this simulation set-up, the alternative variations  
741 of the empirical Bayes method have been excluded and only 50 sets of summary statistics were  
742 simulated for each setting.

743 **S9 Fig. Estimated change in RMSE of significant SNPs at threshold  $5 \times 10^{-4}$  for each method and**  
744 **simulation setting, assuming a quantitative trait, independent SNPs and a skewed effect size**  
745 **distribution.** The estimated change in RMSE of significant SNPs at a threshold of  $5 \times 10^{-4}$  (y-axis),  
746 as defined by Eq (13) in the main text, is plotted for each correction method (x-axis), for each of the  
747 eight simulation settings with a selection coefficient of zero. This figure corresponds to simulation  
748 settings in which it is assumed that the trait of interest is quantitative, SNPs are independent and the  
749 effect sizes follow a skewed distribution. Note that for this simulation set-up, the alternative variations  
750 of the empirical Bayes method have been excluded and only 50 sets of summary statistics were  
751 simulated for each setting.

752 **S10 Fig. Estimated change in RMSE of significant SNPs at threshold  $5 \times 10^{-8}$  for each method**  
753 **and simulation setting, assuming a binary trait, independent SNPs and a normal effect size**  
754 **distribution.** The estimated change in RMSE of significant SNPs at a threshold of  $5 \times 10^{-8}$  (y-axis),  
755 as defined by Eq (13) in the main text, is plotted for each correction method (x-axis), for each of the  
756 eight simulation settings with a selection coefficient of zero. This figure corresponds to simulation  
757 settings in which it is assumed that the trait of interest is binary, SNPs are independent and the effect  
758 sizes follow a normal distribution. Note that for this simulation set-up, the alternative variations of the  
759 empirical Bayes method have been excluded and only 50 sets of summary statistics were simulated  
760 for each setting.

761 **S11 Fig. Estimated change in RMSE of significant SNPs at threshold  $5 \times 10^{-4}$  for each method**  
762 **and simulation setting, assuming a binary trait, independent SNPs and a normal effect size**  
763 **distribution.** The estimated change in RMSE of significant SNPs at a threshold of  $5 \times 10^{-4}$  (y-axis),  
764 as defined by Eq (13) in the main text, is plotted for each correction method (x-axis), for each of the  
765 eight simulation settings with a selection coefficient of zero. This figure corresponds to simulation  
766 settings in which it is assumed that the trait of interest is binary, SNPs are independent and the effect  
767 sizes follow a normal distribution. Note that for this simulation set-up, the alternative variations of the  
768 empirical Bayes method have been excluded and only 50 sets of summary statistics were simulated  
769 for each setting.

770 **S12 Fig. Z-statistics plotted against bias for each real data set.** For each of the six real data sets, z-  
771 statistics (x-axis) are plotted against estimated bias (y-axis). The z-statistic of a SNP is defined as its  
772 estimated effect size divided by the standard error of that estimated effect size while the estimated  
773 bias of a SNP is defined by Eq (15) in the main text. The title of each plot (A)-(F) indicates which  
774 real data set the plot relates to. In each plot, the dark red dashed vertical line represents the z-statistic  
775 corresponding to a p-value of  $5 \times 10^{-8}$  and thus, any points outside these two dark red lines are SNPs  
776 with p-values passing the genome-wide significance threshold of  $5 \times 10^{-8}$ . In a similar manner, the  
777 light red dashed vertical line represents the greater significance threshold of  $5 \times 10^{-4}$ . The dark grey  
778 points highlight SNPs that have p-values less than  $5 \times 10^{-4}$  and have significantly overestimated effect  
779 sizes. The estimated effect size of a SNP has been defined as significantly overestimated according to  
780 Eq (15) in the main text, but in which the true effect size is replaced by the estimated effect size  
781 obtained in the corresponding replication GWAS.

782 **S1 Table. The average number and MSE of significant SNPs at two significance thresholds,  $5 \times$**   
783  **$10^{-8}$  and  $5 \times 10^{-4}$ , with proportions that indicate the extent of *Winner's Curse* bias for each**  
784 **simulation scenario.** S1 Table details an initial exploration of the various simulation scenarios. The  
785 values provided in the table are averages obtained across 100 simulated sets of summary statistics for  
786 each scenario. The top portion of the table shows the values of the parameters which define each  
787 simulation scenario, i.e. sample size, heritability and polygenicity (proportion of effect SNPs). This

788 table corresponds to simulation settings in which a simple correlation structure has been imposed on  
789 the set of SNPs. The number of significant SNPs identified for each scenario at two significance  
790 thresholds,  $5 \times 10^{-8}$  and  $5 \times 10^{-4}$ , is provided, as well as the naive MSE of these significant SNPs, as  
791 defined by Eq (15) in the main manuscript. The table also contains the proportion of significant SNPs  
792 that were seen to have a larger estimated effect size than their true effect size, in terms of absolute  
793 value. The final row for each threshold provides the proportion of significant SNPs that have  
794 significantly overestimated effect sizes. The estimated effect size of a SNP has been defined as  
795 significantly overestimated according to Eq (15) in the main text. This metric gives an indication of  
796 the extent of *Winner's Curse* bias present for each simulation scenario and threshold.

797 **S2 Table. Estimated change in RMSE of significant SNPs at threshold  $5 \times 10^{-8}$  for each method**  
798 **and simulation setting.** S2 Table provides values for the estimated change in RMSE of significant  
799 SNPs, as defined by Eq (14) in the main manuscript, for each *Winner's Curse* correction method and  
800 simulation scenario. This table corresponds to simulation settings in which a simple correlation  
801 structure has been imposed on the set of SNPs. The top portion of the table shows the values of the  
802 parameters which define each simulation scenario, i.e. sample size, heritability and polygenicity  
803 (proportion of effect SNPs). As described in the main manuscript, 100 sets of summary statistics were  
804 simulated for each scenario and the correction methods were applied to each set. Thus, the values  
805 shown in the remaining portion of the table are the average estimated change in RMSE of significant  
806 SNPs due to method implementation across each of these 100 sets. As it is the change in RMSE that  
807 has been computed, it is desirable to obtain a negative change, i.e. the RMSE computed upon  
808 application of the correction method is smaller than that of the naïve approach. Thus, positive values  
809 in the table have been shaded in grey, indicating poor performing methods. The light green shaded  
810 cells highlight the method which, on average, resulted in the greatest reduction in RMSE for each  
811 simulated scenario.

812 **S3 Table. Estimated change in MSE of significant SNPs at threshold  $5 \times 10^{-8}$  for each method**  
813 **and simulation setting.** S3 Table provides values for the estimated change in MSE of significant  
814 SNPs, as defined by Eq (14) in the main manuscript, for each *Winner's Curse* correction method and



815 simulation scenario. This table corresponds to simulation settings in which a simple correlation  
816 structure has been imposed on the set of SNPs. The top portion of the table shows the values of the  
817 parameters which define each simulation scenario, i.e. sample size, heritability and polygenicity  
818 (proportion of effect SNPs). As described in the main manuscript, 100 sets of summary statistics were  
819 simulated for each scenario and the correction methods were applied to each set. Thus, the values  
820 shown in the remaining portion of the table are the average estimated change in MSE of significant  
821 SNPs due to method implementation across each of these 100 sets. As it is the change in MSE that has  
822 been computed, it is desirable to obtain a negative change, i.e. the MSE computed upon application of  
823 the correction method is smaller than that of the naïve approach. Thus, positive values in the table  
824 have been shaded in grey, indicating poor performing methods. The light green shaded cells highlight  
825 the method which, on average, resulted in the greatest reduction in MSE for each simulated scenario.

826 **S4 Table. Estimated relative change in MSE of significant SNPs at threshold  $5 \times 10^{-8}$  for each**  
827 **method and simulation setting, with a simple correlation structure imposed on the set of SNPs.**

828 S4 Table provides values for the estimated relative change in MSE of significant SNPs for each  
829 *Winner's Curse* correction method and simulation scenario. This table corresponds to simulation  
830 settings in which a simple correlation structure has been imposed on the set of SNPs. The top portion  
831 of the table shows the values of the parameters which define each simulation scenario, i.e. sample  
832 size, heritability and polygenicity (proportion of effect SNPs). As described in the main manuscript,  
833 100 sets of summary statistics were simulated for each scenario and the correction methods were  
834 applied to each set. Thus, the values shown in the remaining portion of the table are the average  
835 estimated relative change in MSE of significant SNPs due to method implementation across each of  
836 these 100 sets. As it is the relative change in MSE that has been computed, it is desirable to obtain a  
837 negative change, i.e. the MSE computed upon application of the correction method is smaller than  
838 that of the naïve approach. Thus, positive values in the table have been shaded in grey, indicating poor  
839 performing methods. The light green shaded cells highlight the method which, on average, resulted in  
840 the greatest relative reduction in MSE for each simulated scenario. As the final column contains the  
841 mean of each row, it shows that the bootstrap method has the greatest average estimated relative



842 reduction in MSE. This value of -0.2608 suggests that on average, the bootstrap method improves the  
843 MSE of significant SNPs by  $\approx 26.08\%$ .

844 **S5 Table. Estimated change in RMSE of significant SNPs at threshold  $5 \times 10^{-4}$  for each method**  
845 **and simulation setting.** S5 Table provides values for the estimated change in RMSE of significant  
846 SNPs, as defined by Eq (14) in the main manuscript, for each *Winner's Curse* correction method and  
847 simulation scenario, when a significance threshold of  $5 \times 10^{-4}$  is used. This table corresponds to  
848 simulation settings in which a simple correlation structure has been imposed on the set of SNPs. The  
849 top portion of the table shows the values of the parameters which define each simulation scenario, i.e.  
850 sample size, heritability and polygenicity (proportion of effect SNPs). As described in the main  
851 manuscript, 100 sets of summary statistics were simulated for each scenario and the correction  
852 methods were applied to each set. Thus, the values shown in the remaining portion of the table are the  
853 average estimated change in RMSE of significant SNPs due to method implementation across each of  
854 these 100 sets. As it is the change in RMSE that has been computed, it is desirable to obtain a  
855 negative change, i.e. the RMSE computed upon application of the correction method is smaller than  
856 that of the naïve approach. Thus, positive values in the table have been shaded in grey, indicating poor  
857 performing methods. The light green shaded cells highlight the method which, on average, resulted in  
858 the greatest reduction in RMSE for each simulated scenario.

859 **S6 Table. Estimated change in MSE of significant SNPs at threshold  $5 \times 10^{-4}$  for each method**  
860 **and simulation setting, with a simple correlation structure imposed on the set of SNPs.** S6 Table  
861 provides values for the estimated change in MSE of significant SNPs, as defined by Eq (14) in the  
862 main manuscript, for each *Winner's Curse* correction method and simulation scenario, when a  
863 significance threshold of  $5 \times 10^{-4}$  is used. This table corresponds to simulation settings in which a  
864 simple correlation structure has been imposed on the set of SNPs. The top portion of the table shows  
865 the values of the parameters which define each simulation scenario, i.e. sample size, heritability and  
866 polygenicity (proportion of effect SNPs). As described in the main manuscript, 100 sets of summary  
867 statistics were simulated for each scenario and the correction methods were applied to each set. Thus,  
868 the values shown in the remaining portion of the table are the average estimated change in MSE of

869 significant SNPs due to method implementation across each of these 100 sets. As it is the change in  
870 MSE that has been computed, it is desirable to obtain a negative change, i.e. the MSE computed upon  
871 application of the correction method is smaller than that of the naïve approach. Thus, positive values  
872 in the table have been shaded in grey, indicating poor performing methods. The light green shaded  
873 cells highlight the method which, on average, resulted in the greatest reduction in MSE for each  
874 simulated scenario.

875 **S7 Table. Estimated relative change in MSE of significant SNPs at threshold  $5 \times 10^{-4}$  for each**  
876 **method and simulation setting, with a simple correlation structure imposed on the set of SNPs.**

877 S7 Table provides values for the estimated relative change in MSE of significant SNPs for each  
878 *Winner's Curse* correction method and simulation scenario, when a significance threshold of  $5 \times 10^{-4}$   
879 is used. This table corresponds to simulation settings in which a simple correlation structure has been  
880 imposed on the set of SNPs. The top portion of the table shows the values of the parameters which  
881 define each simulation scenario, i.e. sample size, heritability and polygenicity (proportion of effect  
882 SNPs). As described in the main manuscript, 100 sets of summary statistics were simulated for each  
883 scenario and the correction methods were applied to each set. Thus, the values shown in the remaining  
884 portion of the table are the average estimated relative change in MSE of significant SNPs due to  
885 method implementation across each of these 100 sets. As it is the relative change in MSE that has  
886 been computed, it is desirable to obtain a negative change, i.e. the MSE computed upon application of  
887 the correction method is smaller than that of the naïve approach. Thus, positive values in the table  
888 have been shaded in grey, indicating poor performing methods. The light green shaded cells highlight  
889 the method which, on average, resulted in the greatest relative reduction in MSE for each simulated  
890 scenario. As the final column contains the mean of each row, it shows that the original empirical  
891 Bayes method has the greatest average estimated relative reduction in MSE, when a significance  
892 threshold of  $5 \times 10^{-4}$  is used. This value of -0.3338 suggests that on average, this form of the empirical  
893 Bayes method improves the MSE of significant SNPs by  $\approx 33.38\%$ .

894 **S8 Table. Estimated MSE of significant SNPs at threshold  $5 \times 10^{-4}$  for each method and data set.**

895 S8 Table provides values for the estimated MSE of significant SNPs, as defined by Eq (14) in the

896 main manuscript, using a threshold of  $5 \times 10^{-4}$ , for each *Winner's Curse* correction method and UK  
897 Biobank data set. The first row of values represents the estimated MSE obtained if the unadjusted  
898 estimated effect sizes of the discovery GWAS are used and no correction method has been applied.  
899 This is followed by rows which are representative of the use of different correction methods, i.e. the  
900 conditional likelihood based methods, the empirical Bayes method and its variations, the proposed  
901 bootstrap method and FIQT, respectively. As it is desirable to obtain lower estimated MSE values  
902 upon application of a method, values which are greater than their corresponding naïve value have  
903 been shaded in grey. The light green shaded cells highlight the method which resulted in the lowest  
904 estimated MSE value for each data set.

905 **S1 File. Text Supplement.** This file contains a more detailed description of the various proposed  
906 modifications to the empirical Bayes method, the simulation process and the evaluation of method  
907 performance using simulated data sets of independent SNPs.